

Precision	Number of Bits	Sign Bit Range	Number of Bits	Exponent Bit Range	Number of Bits	Mantissa Bit Range	Bias
Single	1	[31]	8	[30-23]	23	[22-0]	127
Double	1	[63]	11	[62-52]	52	[51-0]	1023

Figure 1A

Floating Point Numbers  
Under IEEE Standard 754

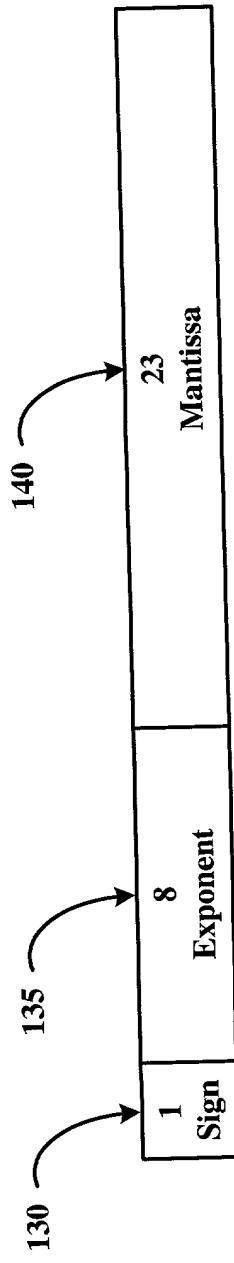
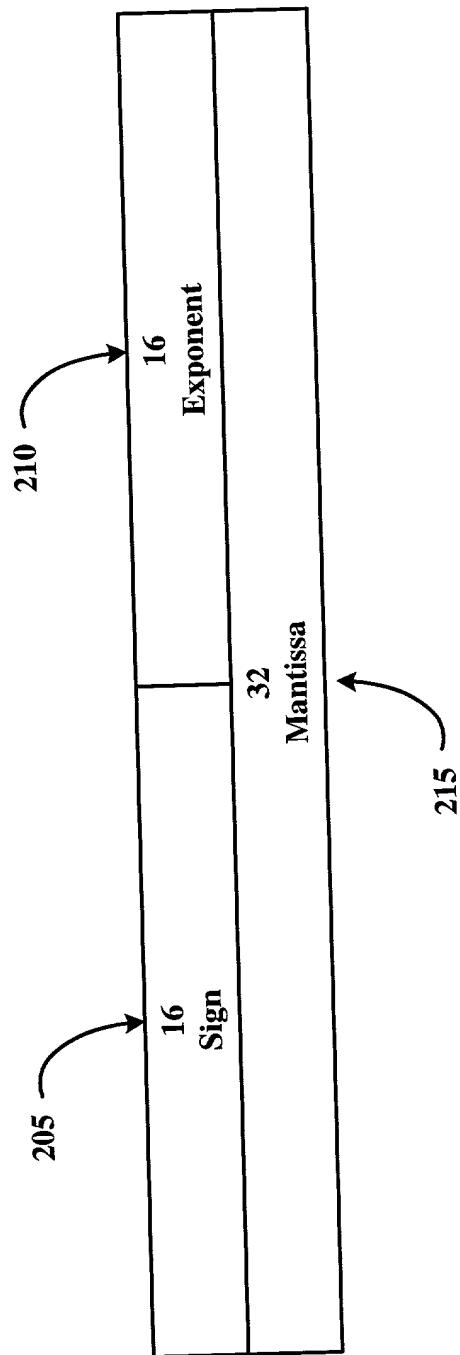


Figure 1B

Floating Point  
Single Precision Number



Unpacked  
Floating Point Number

Figure 2

```

305
    FLO_type
    add (FLO_type arg_a, FLO_type arg_b)
    {
        fp_number_type a;
        fp_number_type b;
        fp_number_type tmp;
        fp_number_type *res;
        FLO_union_type au, bu;

        au.value = arg_a;
        bu.value = arg_b;

        unpack_d (&au, &a);
        unpack_d (&bu, &b);

        315    → res = _fpadd_parts (&a, &b, &tmp);
        320    → unpack_d (&au, &a);
        320    → unpack_d (&bu, &b);
        330    → b.sign ^= 1;
        335    → res = _fpadd_parts (&a, &b, &tmp);
        345    → return pack_d (res);
    }

```

```

310
    FLO_type
    sub (FLO_type arg_a, FLO_type arg_b)
    {
        fp_number_type a;
        fp_number_type b;
        fp_number_type tmp;
        fp_number_type *res;
        FLO_union_type au, bu;

        au.value = arg_a;
        bu.value = arg_b;

        unpack_d (&au, &a);
        unpack_d (&bu, &b);

        315    → res = _fpadd_parts (&a, &b, &tmp);
        320    → unpack_d (&au, &a);
        320    → unpack_d (&bu, &b);
        330    → b.sign ^= 1;
        335    → res = _fpadd_parts (&a, &b, &tmp);
        345    → return pack_d (res);
    }

```

Figure 3A

Conventional Addition Routine

Figure 3B

Conventional Subtraction Routine

<b>405</b>	$T = a \times b \times c$
<b>410</b>	$T_0 = a \times b$
<b>415</b>	$T_1 = T_0 \times c$
<b>420</b>	$T_2 = \text{unpack}(a)$
<b>425</b>	$T_3 = \text{unpack}(b)$
<b>430</b>	$T_4 = \text{unpack\_mult}(T_2, T_3)$
<b>435</b>	$T_0 = \text{pack}(T_4)$
<b>440</b>	$T_5 = \text{unpack}(T_0)$
<b>445</b>	$T_6 = \text{unpack}(c)$
<b>450</b>	$T_7 = \text{unpack\_mult}(T_5, T_6)$
<b>455</b>	$T_1 = \text{pack}(T_7)$

Figure 4

Calculation Using Conventional  
Floating Point Emulation

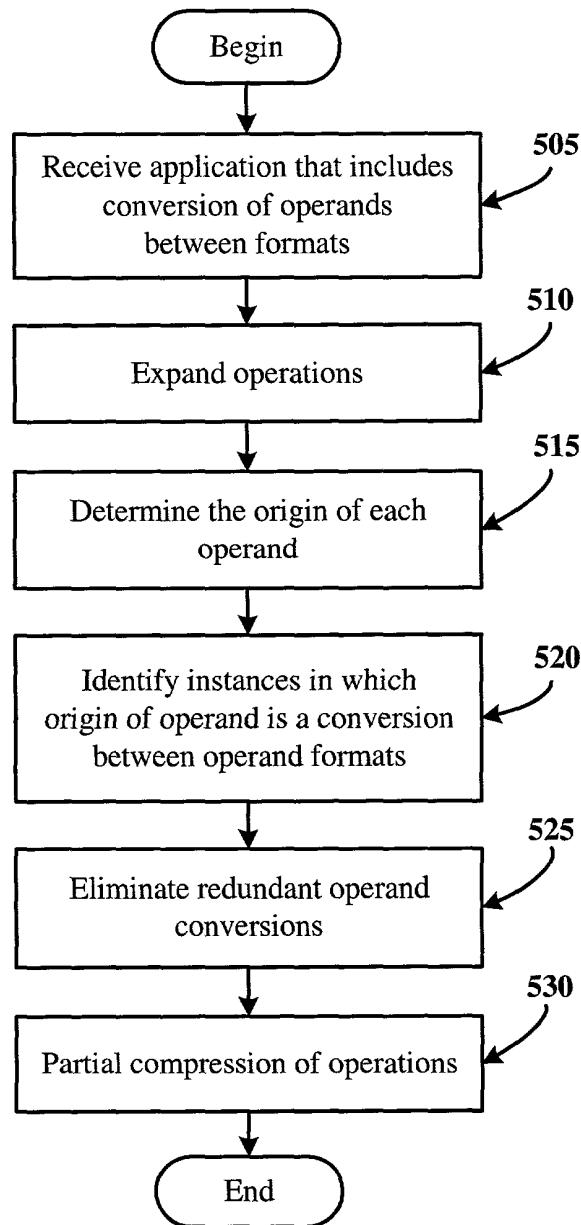


Figure 5

Operand Conversion Optimization

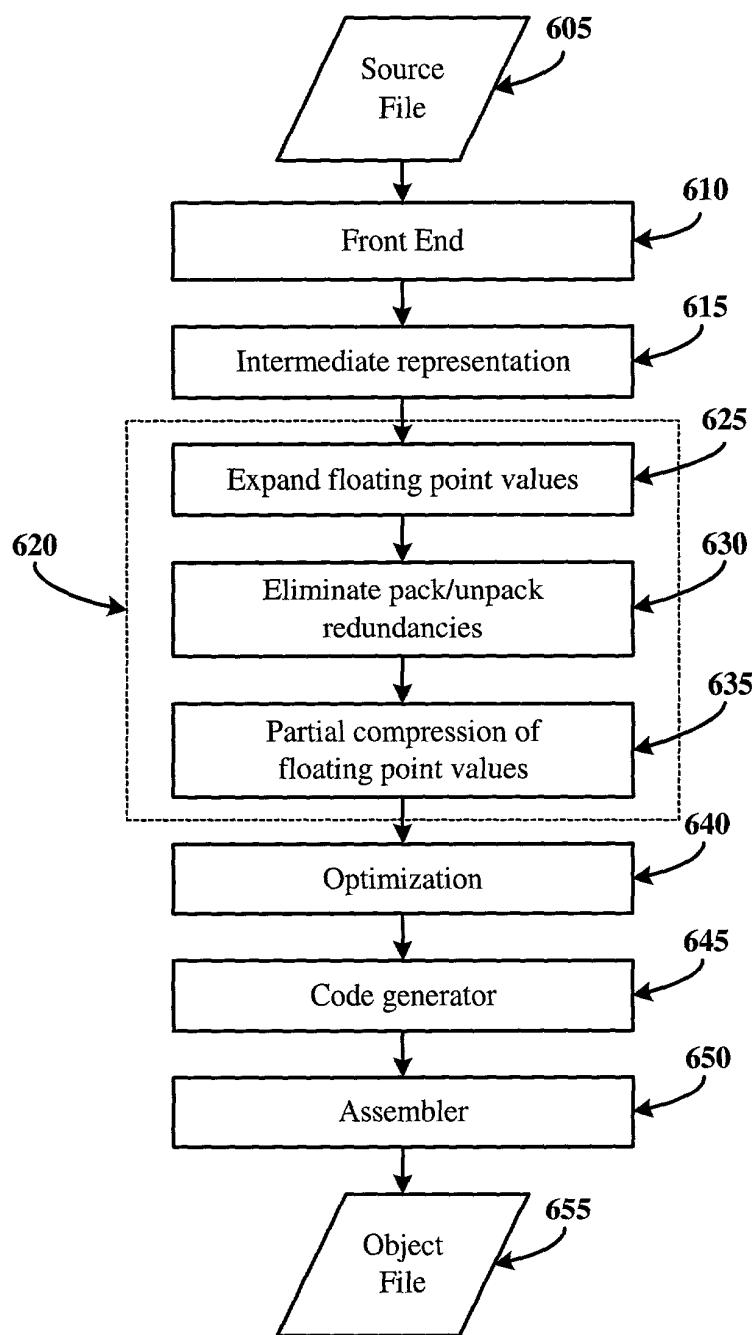


Figure 6

Exemplary System Including  
Operand Conversion Optimization

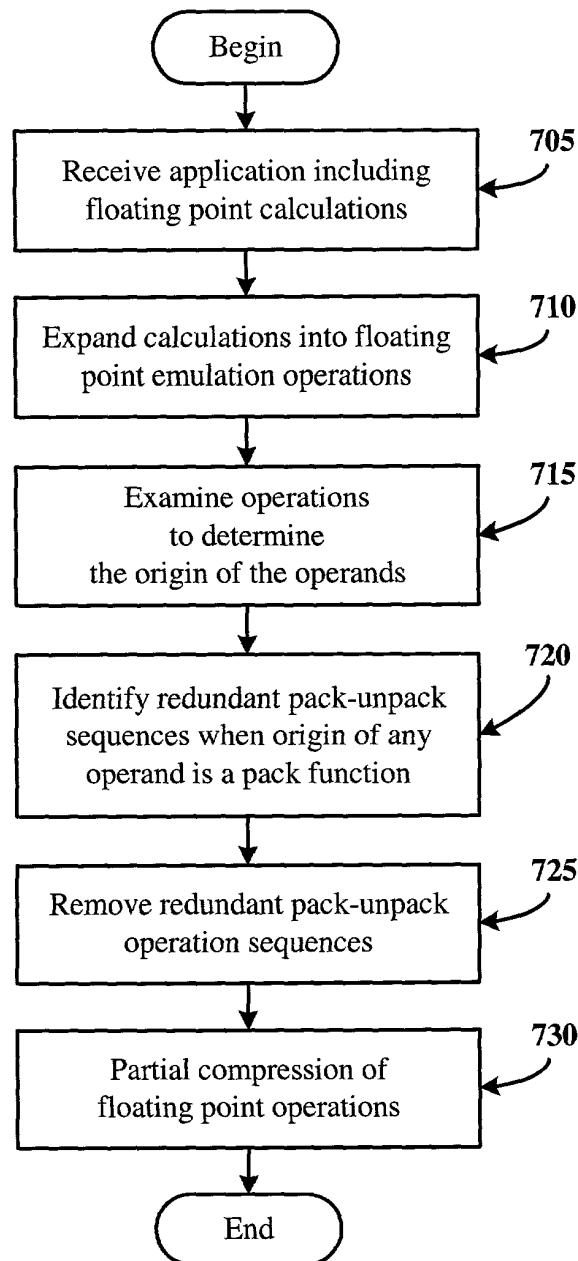


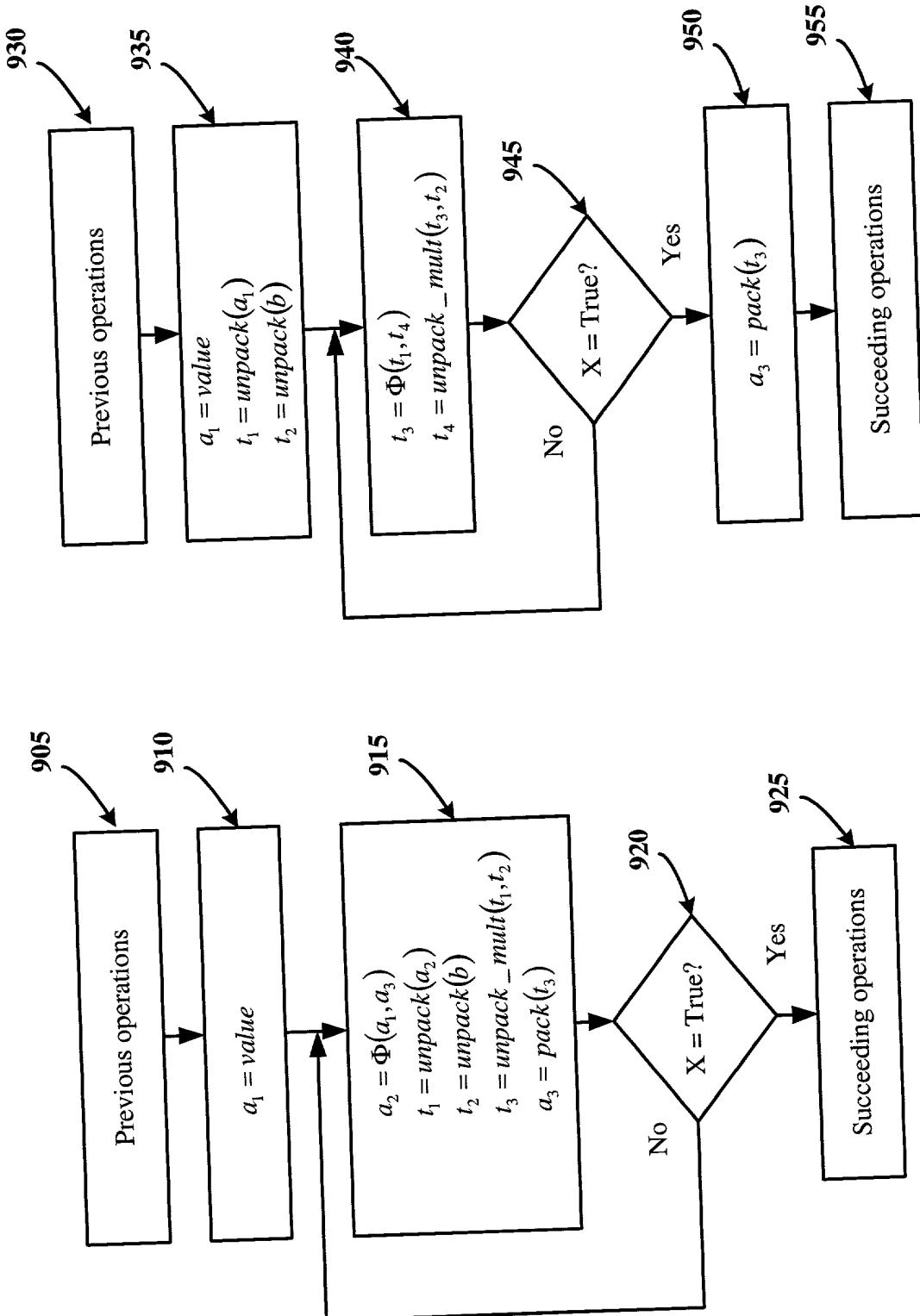
Figure 7

Optimization of  
Floating Point Emulation

<b>805</b>	$T = a \times b \times c$
<b>810</b>	$T_0 = a \times b$
<b>815</b>	$T_1 = T_0 \times c$
<b>820</b>	$T_2 = \text{unpack}(a)$
<b>825</b>	$T_3 = \text{unpack}(b)$
<b>830</b>	$T_4 = \text{unpack\_mult}(T_2, T_3)$
<b>835</b>	$T_5 = \text{unpack}(c)$
<b>840</b>	$T_6 = \text{unpack\_mult}(T_4, T_5)$
<b>845</b>	$T_1 = \text{pack}(T_6)$

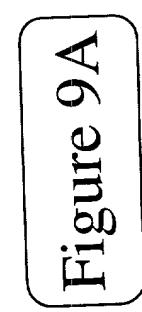
Figure 8

Calculation Using  
Emulation Optimization



Loop with  
Partial Redundancy

Figure 9A  
Optimized Loop



Loop with  
Partial Redundancy

Figure 9B  
Optimized Loop